

# DEEPLIFT

Monday, 20th of May, 2019

**Backpropagation Method** : A method for decomposing the output prediction of a neural network on a specific input by **backpropagating the contributions** of all neurons in the network to every feature of the input

**Reference point** : the difference in output from some 'reference' output in terms of the difference of the input from some 'reference' input

let  $t$  represent some target output neuron of interest

Let  $x_1, x_2, \dots, x_n$  represent some neurons in some intermediate layer or set of layers that are necessary and sufficient to compute  $t$ .

Let  $t_0$  represent the reference activation of  $t$ .

$\Delta t$  : difference-from-reference, that is  $\Delta t = t - t_0$

**Contribution score** : Amount of difference-from- reference in  $t$  that is attributed to or ‘blamed’ on the difference-from-reference of  $x$

**Summation of  $C_{\Delta x i \Delta t} = \Delta t$**

**Choose Reference point** : As a guiding principle, we can ask ourselves “what am I interested in measuring differences against?”.

Example: For MNIST, we use a reference input of all zeros as this is the background of the images.

**Multipliers** :  $m_{\Delta x \Delta t} = C_{\Delta x \Delta t} / \Delta x$

**Chain Rule for Multipliers** :  $m_{\Delta x i \Delta z} = \text{summation over } j (m_{\Delta x i \Delta y j} * m_{\Delta y j \Delta z})$

**Defining Reference** : references for all neurons can be found by choosing a reference input and propagating activations through the net.

# Rules for assigning contribution scores

## THE LINEAR RULE :

Dense and Convolutional layers (excluding nonlinearities). Let  $y$  be a linear function of its inputs  $x_i$  such that  $y = b + \textit{summation over } i (w_i x_i)$ . We have  $\Delta y = \textit{summation over } i (w_i \Delta x_i)$ . We define the positive and negative parts of  $\Delta y$  as:

$$\Delta y^+ = \textit{summation over } i (1_{\{w_i \Delta x_i > 0\}} w_i \Delta x_i)$$

$$\Delta y^- = \textit{summation over } i (1_{\{w_i \Delta x_i < 0\}} w_i \Delta x_i)$$

$$C_{\Delta x+ i \Delta y+} = 1\{w_i \Delta x_i > 0\} w_i \Delta x+ i$$

$$C_{\Delta x- i \Delta y+} = 1\{w_i \Delta x_i > 0\} w_i \Delta x- i$$

$$C_{\Delta x+ i \Delta y-} = 1\{w_i \Delta x_i < 0\} w_i \Delta x+ i$$

$$C_{\Delta x- i \Delta y-} = 1\{w_i \Delta x_i < 0\} w_i \Delta x- i$$

$$m_{\Delta x+ i \Delta y+} = m_{\Delta x- i \Delta y+} = 1\{w_i \Delta x_i > 0\} w_i$$

$$m_{\Delta x+ i \Delta y-} = m_{\Delta x- i \Delta y-} = 1\{w_i \Delta x_i < 0\} w_i$$

**THE RESCALE RULE :** This rule applies to nonlinear transformations that take a single input , such as the ReLU, tanh or sigmoid operations.

$$\Delta y+ = \Delta y / \Delta x * \Delta x+ = C_{\Delta x+} \Delta y+$$

$$\Delta y- = \Delta y / \Delta x * \Delta x- = C_{\Delta x-} \Delta y-$$

Based on this, we get:

$$m_{\Delta x+} \Delta y+ = m_{\Delta x-} \Delta y- = m_{\Delta x} \Delta y = \Delta y / \Delta x$$

Can give misleading results, as it calculates the contribution of  $\Delta x^+$  *in presence of*  $\Delta x^-$  - which sometimes lead to **zero contributions of both  $\Delta x^+$  and  $\Delta x^-$**  .

**Where we might prefer to use the Rescale rule?**

Specifically, consider a thresholded ReLU where  $\Delta y > 0$  *iff*  $\Delta x \geq b$ . If  $\Delta x < b$  merely indicates noise, we would want to assign contributions of 0 to both  $\Delta x^+$  and  $\Delta x^-$  (as done by the Rescale rule) to mitigate the noise.

## THE REVEAL CANCEL RULE:

By considering the *impact of the positive terms in the absence of negative terms, and the impact of negative terms in the absence of positive terms*, we alleviate some of the issues that arise from positive and negative terms canceling each other out.

$$\Delta y^+ = \frac{1}{2} * (f(x_0 + \Delta x^+) - f(x_0)) + \frac{1}{2} * (f(x_0 + \Delta x^- + \Delta x^+) - f(x_0 + \Delta x^-))$$

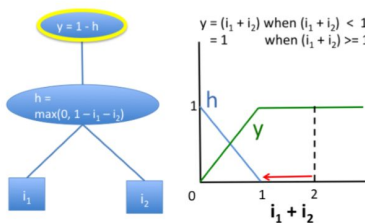
$$\Delta y^- = \frac{1}{2} * (f(x_0 + \Delta x^-) - f(x_0)) + \frac{1}{2} * (f(x_0 + \Delta x^+ + \Delta x^-) - f(x_0 + \Delta x^+))$$

$$m_{\Delta x^+ \Delta y^+} = C_{\Delta x^+ y^+} / \Delta x^+ = \Delta y^+ / \Delta x^+; \quad m_{\Delta x^- \Delta y^-} = \Delta y^- / \Delta x^-$$

# Problems with other back-propagation methods using gradients (like Layerwise Relevance Propagation)

1. **Saturation problem** : When value of the output function becomes constant or it saturates i.e gradient becomes zero, after a certain input value, the importance of the input features lying after that saturation point also becomes zero which is a misleading result.

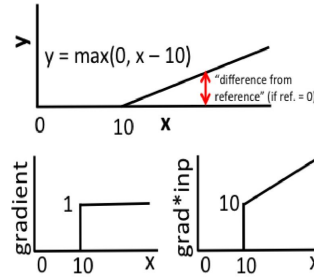
However, the contribution scores calculated by the DEELIFT method can be non zero even if gradient at a point is zero.



**Figure 1. Perturbation-based approaches and gradient-based approaches fail to model saturation.** Illustrated is a simple network exhibiting saturation in the signal from its inputs. At the point where  $i_1 = 1$  and  $i_2 = 1$ , perturbing either  $i_1$  or  $i_2$  to 0 will not produce a change in the output. Note that the gradient of the output w.r.t the inputs is also zero when  $i_1 + i_2 > 1$ .

**2. Thresholding problem :** Sometimes there is a sudden change in the gradients such as in a following linear graph, which gives contribution scores with larger difference even in the infinitesimally small change in the input.

However, DEELIFT method doesn't arise such problems due to the concept of difference-from-reference. It gives continuous values in such regions.



*Figure 2. Discontinuous gradients can produce misleading importance scores.* Response of a single rectified linear unit with a bias of  $-10$ . Both gradient and gradient $\times$ input have a discontinuity at  $x = 10$ ; at  $x = 10 + \epsilon$ , gradient $\times$ input assigns a contribution of  $10 + \epsilon$  to  $x$  and  $-10$  to the bias term ( $\epsilon$  is a small positive number). When  $x < 10$ , contributions on  $x$  and the bias term are both 0. By contrast, the difference-from-reference (red arrow, top figure) gives a continuous increase in the contribution score.

Github repository : <https://github.com/kundajelab/deeplift>



## Work-Flow :

1. **Deeplift** : “Learning Important Features Through Propagating Activation Differences” by Shrikumar, Greenside & Kundaje
2. **Internal Influence based approach** : K. Leino, L. Li, S. Sen, A. Datta, and M. Fredrikson. Influence-Directed Explanations for Deep Convolutional Networks. In CoRR abs/1802.03788, 2018.
  - a. Comparison between both
  - b. Comparison with Layerwise Relevance Propagation

### 3. Self Explaining Models :

Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV)

David Alvarez-Melis, Tommi S. Jaakkola, Towards Robust Interpretability with Self-Explaining Neural Networks

# Requirements:

GPU server access (CUDA) - LAPTOP

NVIDIA driver + CUDA + tensorflow + Opencv + LINUX + python - SYSTEM